



TITLE:

Note on robust model selection by density power divergence in a contaminated regression model (Statistical Information in Inference and Its Related Topics)

AUTHOR(S):

矢田, 和善; 青嶋, 誠; 小林, 裕子

---

CITATION:

矢田, 和善 ...[et al]. Note on robust model selection by density power divergence in a contaminated regression model (Statistical Information in Inference and Its Related Topics). 数理解析研究所講究録 2011, 1758: 150-159

ISSUE DATE:

2011-08

URL:

<http://hdl.handle.net/2433/171314>

RIGHT:

## Note on robust model selection by density power divergence in a contaminated regression model

筑波大学・数学系 矢田 和善 (Kazuyoshi Yata)  
Institute of Mathematics  
University of Tsukuba

筑波大学・数学系 青嶋 誠 (Makoto Aoshima)  
Institute of Mathematics  
University of Tsukuba

筑波大学大学院・数理物質科学研究科 小林 裕子 (Yuko Kobayashi)  
Graduate School of Pure and Applied Sciences  
University of Tsukuba

### 1. はじめに

パラメトリックなモデルの推定において、ダイバージェンス最小化による推定法は長い歴史をもつ。最もよく知られるのは、Kullback-Leibler 情報量 (K-L 情報量) 最小化に基づく最尤推定であろう。標本数が十分ではないとき、最尤推定量 (MLE) は異常値の影響を受けやすいという欠点をもつ。一方, Hjort (1994), Scott (2001) は,  $L_2$  距離最小化に基づく  $L_2$  推定量を提案した。  $L_2$  推定量は異常値に対して頑健ではあるものの, 漸近有効ではないという欠点をもつ。これらに対して, Basu et al. (1998) は density power divergence (DPD) を提案した。DPD は, チューニングパラメータ  $\alpha (\geq 0)$  の値を変えることによって頑健性と漸近有効性のトレードオフを調節した推定を考えることができ,  $\alpha = 0$  のときに K-L 情報量と一致し,  $\alpha = 1$  のときに  $L_2$  距離と一致する。Jones et al. (2001), Basu et al. (2006) は, DPD 最小化によって与えられる推定量の漸近的性質を研究した。関連する研究として, Fujisawa and Eguchi (2006) は,  $\beta$ -ダイバージェンスという DPD と同等なダイバージェンスに基づいて, 1 次元混合正規分布のパラメータを推定するためのアルゴリズムを提案した。また, Fujisawa and Eguchi (2008) では,  $\gamma$ -ダイバージェンスという頑健な推定を考えた。

DPD 最小化による推定は, 漸近有効ではあるが頑健でない MLE と, 頑健ではあるが漸近有効でない  $L_2$  推定との間の, 架け橋として考えることができる。頑健性と漸近有効性を調節するためには,  $\alpha$  の値を適切に決めることが重要である。 $\alpha$  の選択法について, Hong and Kim (2001), Warwick and Jones (2005), Warwick (2005), Durio and Isaia (2011) 等の研究があるが, モデル選択の立場からは, 使い勝手がいいとは必ずしも言い難い。最近, Kobayashi et al. (2011) は, 異常値混入モデルにおける潜在分布のクラス数とモデルを有効に推定するための  $\alpha$  の選択法

を提案した。そこでは、異常値混入モデルにおける DPD ( $\beta$ -ダイバージェンス) を漸近的に評価し、潜在分布のモデル選択のための新しい情報量規準を与えた。

本稿では、Kobayashi et al. (2011) が与えた DPD の漸近理論に基づいて、回帰モデルにおける頑健なモデル選択のための、新しい情報量規準を提案する。提案する情報量規準が、異常値に対して頑健なモデル選択を保証することを、理論的かつ数値的に考察する。

## 2. 異常値に対して頑健な推定量

$d$  次元データ  $\mathbf{x}$  を生成する真の分布の密度関数を  $g(\mathbf{x})$  とし、その分布関数を  $G(\mathbf{x})$  とする。  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^T$  はあるパラメータ数  $p$  をもつパラメータベクトルとする。未知の  $g(\mathbf{x})$  を近似するモデルを  $f(\mathbf{x}|\boldsymbol{\theta})$  とし、これを異常値に対して頑健に構築したい。そのために、  $f(\mathbf{x}|\boldsymbol{\theta})$  の  $g(\mathbf{x})$  に対する近さを

$$D_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta})) = \begin{cases} \frac{1}{\alpha} \int_{\mathbf{R}^d} g(\mathbf{x})^{1+\alpha} d\mathbf{x} - \left(1 + \frac{1}{\alpha}\right) \int_{\mathbf{R}^d} f(\mathbf{x}|\boldsymbol{\theta})^\alpha g(\mathbf{x}) d\mathbf{x} \\ + \int_{\mathbf{R}^d} f(\mathbf{x}|\boldsymbol{\theta})^{1+\alpha} d\mathbf{x} & (\alpha > 0), \\ \int_{\mathbf{R}^d} (\log g(\mathbf{x}) - \log f(\mathbf{x}|\boldsymbol{\theta})) g(\mathbf{x}) d\mathbf{x} & (\alpha = 0) \end{cases} \quad (2.1)$$

で測る。  $D_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta}))$  は density power divergence(DPD) といい、次の性質をもつ。

$$(i) D_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta})) \geq 0$$

$$(ii) D_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta})) = 0 \iff g(\mathbf{x}) = f(\mathbf{x}|\boldsymbol{\theta}), \forall \mathbf{x} \in \mathbf{R}^d$$

$D_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta}))$  が小さいほど、モデル  $f(\mathbf{x}|\boldsymbol{\theta})$  は真の  $g(\mathbf{x})$  に近いと考えられる。これ以降は、  $\alpha > 0$  の場合を扱う。いま、

$$C_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta})) = \int_{\mathbf{R}^d} f(\mathbf{x}|\boldsymbol{\theta})^{1+\alpha} d\mathbf{x} - \left(1 + \frac{1}{\alpha}\right) \int_{\mathbf{R}^d} f(\mathbf{x}|\boldsymbol{\theta})^\alpha g(\mathbf{x}) d\mathbf{x} \quad (2.2)$$

とおく。これを用いて

$$D_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta})) = C_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta})) - C_\alpha(g(\mathbf{x}); g(\mathbf{x})) \quad (2.3)$$

と表せるので、  $D_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta}))$  を最小にするパラメータは、

$$\boldsymbol{\theta}_\alpha = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} C_\alpha(g(\mathbf{x}); f(\mathbf{x}|\boldsymbol{\theta}))$$

で与えられる。ただし、  $\Theta$  は  $\boldsymbol{\theta}$  の開空間である。(2.2) 式における  $g(\mathbf{x})$  の分布関数  $G(\mathbf{x})$  を経験分布関数  $\hat{G}(\mathbf{x})$  で置き換え、

$$\ell_\alpha(\boldsymbol{\theta}; f(\mathbf{x}|\boldsymbol{\theta})) = \int_{\mathbf{R}^d} f(\mathbf{x}|\boldsymbol{\theta})^{1+\alpha} d\mathbf{x} - \left(1 + \frac{1}{\alpha}\right) n^{-1} \sum_{i=1}^n f(\mathbf{x}_i|\boldsymbol{\theta})^\alpha \quad (2.4)$$

とおく. そのとき,

$$\hat{\theta}_\alpha = \operatorname{argmin}_{\theta \in \Theta} \ell_\alpha(\theta; f(\mathbf{x}|\theta))$$

を用いて構築したモデル  $f(\mathbf{x}|\hat{\theta}_\alpha)$  は, データ数  $n$  が大きいときに漸近的に最適なモデルと考えられる.  $\hat{\theta}_\alpha$  は最小 DPD 推定量 (MDPDE) とよばれ,  $\alpha = 0$  のときに最尤推定量 (MLE) と一致し,  $\alpha = 1$  のときに Scott (2001) が提案した  $L_2$  推定量と一致する.

いま,  $I_\alpha(\theta)$ ,  $J_\alpha(\theta)$  を

$$I_\alpha(\theta) = \int_{\mathbf{R}^d} \frac{\partial \ell_\alpha(\theta; f(\mathbf{x}|\theta))}{\partial \theta} \frac{\partial \ell_\alpha(\theta; f(\mathbf{x}|\theta))}{\partial \theta^T} dG(\mathbf{x}),$$

$$J_\alpha(\theta) = \int_{\mathbf{R}^d} \frac{\partial^2 \ell_\alpha(\theta; f(\mathbf{x}|\theta))}{\partial \theta \partial \theta^T} dG(\mathbf{x})$$

と定義し,  $f(\mathbf{x}|\theta)$  について次の正則条件を仮定する:

- (A1)  $f(\mathbf{x}|\theta)$  の台  $A = \{\mathbf{x} | f(\mathbf{x}|\theta) > 0\}$  は  $\theta$  に無関係である.
- (A2)  $\int_{\mathbf{R}^d} f(\mathbf{x}|\theta)^{1+\alpha} d\mathbf{x} < \infty$  であり, すべての  $\mathbf{x}$  で  $f(\mathbf{x}|\theta)$  は  $\theta$  に関して 3 階連続微分可能である. また,  $\theta$  に関する微分と  $\mathbf{x}$  に関する積分の順序交換可能である.
- (A3)  $J_\alpha(\theta_\alpha) > O$  であり, すべての成分は有界である.
- (A4) すべての  $\theta \in \Theta$  について

$$\left| \frac{\partial^3 \ell_\alpha(\theta; f(\mathbf{x}|\theta))}{\partial \theta_i \partial \theta_j \partial \theta_l} \right| \leq M_{ijl}(\mathbf{x}), \quad E_G[M_{ijl}(\mathbf{x})] < \infty; \quad i, j, l = 1, \dots, p$$

となる  $M_{ijl}(\mathbf{x})$  が存在する.

そのとき, MDPDE は次の漸近正規性をもつ.

$$\sqrt{n}(\hat{\theta}_\alpha - \theta_\alpha) \xrightarrow{d} N_p(0, J_\alpha(\theta_\alpha)^{-1} I_\alpha(\theta_\alpha) J_\alpha(\theta_\alpha)^{-1}) \quad \text{as } n \rightarrow \infty. \quad (2.5)$$

次に, MDPDE の異常値に対する頑健性を, シミュレーション実験で検証する. 真の分布に, 次のような異常値が混入した分布を考えた.

$$g(x) = 0.97\phi(x|0, 1) + 0.03\psi(x|4 < x < 8)$$

ここで,  $\phi(x|0, 1)$  は潜在分布を表し, 標準正規分布  $N(0, 1)$  を仮定した. 一方,  $\psi(x|4 < x < 8)$  は異常値の分布を表し, 区間  $(4, 8)$  の一様分布  $U(4, 8)$  を仮定した. つまり, 真の分布には, 3% の確率で異常値が混入した状況を考えた. 真の分布の平均と分散は,  $E(x) = 0.180$ ,  $V(x) = 2.058$  である.  $g(x)$  から独立にデータを 100 個発生させ, 真の分布の平均と分散を MLE ( $\alpha = 0$ ) と MDPDE ( $\alpha = 0.1, 0.3$ ,

0.9)で推定した。これを独立に1000回繰り返し、推定値の平均とその分散（括弧内）を纏めたものが表1である。MLEは、異常値の影響を忠実に反映して、異常値が混在した真の分布の平均と分散を推定した。それに対して、MDPDEは、異常値の影響を抑えて潜在分布の平均と分散を推定する様子が見て取れた。しかしながら、 $\alpha$ の値が小さいとき、MDPDEは異常値の影響を十分に抑えられているわけではなく、また、 $\alpha$ の値が大き過ぎると、推定値の分散が大きくなり推定が不安定になることも見て取れた。したがって、MDPDEの使用には、適切な $\alpha$ の値を選択することが重要であるといえる。

表1. MLE ( $\alpha = 0$ ) と MDPDE ( $\alpha=0.1, 0.3, 0.9$ ) による  
平均と分散の推定値とその分散（括弧内）

| $\theta \setminus \alpha$ | 0                | 0.1              | 0.3               | 0.9               |
|---------------------------|------------------|------------------|-------------------|-------------------|
| $\mu = 0$                 | 0.197<br>(0.024) | 0.059<br>(0.016) | -0.014<br>(0.012) | -0.019<br>(0.015) |
| $\sigma^2 = 1$            | 2.191<br>(0.480) | 1.418<br>(0.199) | 0.990<br>(0.028)  | 0.986<br>(0.050)  |

### 3. 回帰モデルにおける MDPDE

いま、目的変数を  $y_i$ 、 $d+1$ 次元の説明変数ベクトルを  $\mathbf{x}_i = (1, x_{i1}, \dots, x_{id})^T$  とし、独立な  $n$  個のデータ  $\{(y_i, \mathbf{x}_i); i = 1, \dots, n\}$  を観測したとする。線形回帰モデル

$$y_i = \boldsymbol{\beta}^T \mathbf{x}_i + \varepsilon_i \quad (3.1)$$

を考える。ここで、 $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_d)^T$  は  $d+1$ 次元の回帰係数ベクトルで、 $\varepsilon_i$ ,  $i = 1, \dots, n$  は独立同分布に従う誤差変数である。誤差変数  $\varepsilon$  の密度関数  $g(\varepsilon)$  は、ある開領域  $\mathbf{R}_0$  に確率  $\tau \in (0, 1)$  で異常値が混入した次の分布を仮定する。

$$g(\varepsilon) = (1 - \tau)\phi(\varepsilon|0, \sigma^2) + \tau\psi(\varepsilon) \quad (3.2)$$

ここで、 $\phi(\varepsilon|0, \sigma^2)$  は潜在分布となる正規分布  $N(\varepsilon|0, \sigma^2)$  の密度関数を表す。一方、 $\psi(\cdot)$  は異常値の密度関数を表し、

$$\int_{\mathbf{R}_0} \psi(\varepsilon) d\varepsilon = 1, \quad \sup_{\varepsilon \in \mathbf{R}_0} \psi(\varepsilon) < \infty \quad (3.3)$$

を満たすとする。異常値が混入する開領域  $\mathbf{R}_0$  における潜在分布の確率

$$\int_{\mathbf{R}_0} \phi(\varepsilon|0, \sigma^2) d\varepsilon (= \delta > 0, \text{ say})$$

は十分小さいと仮定する。我々は、異常値に頑健な  $\boldsymbol{\beta}$  と  $\sigma^2$  の推定量を MDPDE を用いて構築する。なお、Durio and Isaia (2011) も回帰モデルにおいて MDPDE を考え、 $\alpha$  の選択法を議論している。

モデル  $f(y_i|x_i; \theta)$  は, パラメータ  $\theta = (\beta^T, \sigma^2)^T$  をもつ正規分布  $N(y_i|x_i; \beta^T x_i, \sigma^2)$

$$f(y_i|x_i; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}\right) \quad (3.4)$$

を仮定する. このとき,  $\alpha > 0$  において

$$\begin{aligned} \ell_\alpha(\theta; f(y|x; \theta)) &= n^{-1} \sum_{i=1}^n \int_{\mathcal{R}} f(y_i|x_i; \theta)^{1+\alpha} dy_i - \left(1 + \frac{1}{\alpha}\right) n^{-1} \sum_{i=1}^n f(y_i|x_i; \theta)^\alpha \\ &= \frac{1}{\sqrt{(2\pi\sigma^2)^\alpha(1+\alpha)}} - \left(1 + \frac{1}{\alpha}\right) n^{-1} \sum_{i=1}^n f(y_i|x_i; \theta)^\alpha \end{aligned} \quad (3.5)$$

と定義して, MDPDE を

$$\hat{\theta}_\alpha = \underset{\theta \in \Theta}{\operatorname{argmin}} \ell_\alpha(\theta; f(y|x; \theta))$$

で求める. なお,  $\alpha = 0$  においては MLE と一致し,  $\mathbf{y} = (y_1, \dots, y_n)^T$ ,  $\mathbf{X} = (x_1, \dots, x_n)$  に対して  $\hat{\beta} = (\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{y}$ ,  $\hat{\sigma}^2 = \sum_{i=1}^n (y_i - \hat{\beta}^T x_i)^2/n$  とおけば,  $\hat{\theta}_0 = (\hat{\beta}^T, \hat{\sigma}^2)^T$  となる.

表 2. MLE ( $\alpha = 0$ ) と MDPDE ( $\alpha=0.2, 0.4, 0.6$ ) による  $\beta$  と  $\sigma^2$  の推定値とその分散 (括弧内)

| $\alpha \setminus \theta$ | $\beta_0 = 1.0$  | $\beta_1 = -1.0$  | $\beta_2 = -0.5$  | $\beta_3 = 0.5$  | $\beta_4 = 1.0$  | $\sigma^2 = 1$   |
|---------------------------|------------------|-------------------|-------------------|------------------|------------------|------------------|
| 0                         | 1.327<br>(0.623) | -1.029<br>(0.152) | -0.467<br>(0.150) | 0.489<br>(0.145) | 0.989<br>(0.131) | 2.545<br>(0.946) |
| 0.2                       | 1.063<br>(0.302) | -1.014<br>(0.069) | -0.483<br>(0.065) | 0.488<br>(0.073) | 0.983<br>(0.058) | 1.104<br>(0.197) |
| 0.4                       | 1.043<br>(0.312) | -1.015<br>(0.071) | -0.485<br>(0.066) | 0.485<br>(0.075) | 0.977<br>(0.061) | 0.919<br>(0.051) |
| 0.6                       | 1.045<br>(0.354) | -1.015<br>(0.080) | -0.485<br>(0.076) | 0.484<br>(0.084) | 0.975<br>(0.070) | 0.888<br>(0.057) |

表 2 は, 異常値に対する MDPDE の頑健性をシミュレーション実験で検証したものである. ここでは,  $d = 4$  の線形回帰モデル

$$y_i = \beta_0 + \sum_{j=1}^4 \beta_j x_{ij} + \varepsilon_i$$

を考え,  $(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4) = (1, -1, -0.5, 0.5, 1)$  と固定し,  $x_{ij}, j = 1, \dots, 4$  は一様分布  $U(0, 2)$  から独立に発生させた. 誤差変数  $\varepsilon$  は, 異常値の混入した以下の分布を考えた.

$$g(\varepsilon) = 0.95\phi(\varepsilon|0, 1) + 0.05\psi(\varepsilon|4 < \varepsilon < 8)$$

ここで、潜在分布にあたる  $\phi(\varepsilon|0, 1)$  は標準正規分布  $N(0, 1)$  とし、異常値の分布にあたる  $\psi(\varepsilon|4 < \varepsilon < 8)$  は一様分布  $U(4, 8)$  とした。  $(y_i, \mathbf{x}_i)$  なるデータを独立に 60 個発生させ、回帰係数  $\beta = (1, -1, -0.5, 0.5, 1)^T$  と潜在分布の分散  $\sigma^2 = 1$  を、MLE ( $\alpha = 0$ ) と MDPDE ( $\alpha=0.2, 0.4, 0.6$ ) で推定した。この推定を独立に 1000 回繰り返し、推定値の平均とその分散 (括弧内) を計算した。切片  $\beta_0$  と分散  $\sigma^2$  の推定において、MLE は異常値の影響を顕著に受けた。一方で、MDPDE は異常値の影響を受けることなく、概ね全てのパラメータを良好に推定している様子が見て取れた。

#### 4. 回帰モデルにおける頑健なモデル選択

線形回帰モデル (3.1) において、 $d$  個の説明変数のうちある  $k_*$  ( $\leq d$ ) 個からなる  $\mathbf{x}_{i(*)} = (1, x_{i(1*)}, \dots, x_{i(k_*)})^T$  によって構築される

$$y_i = \beta_{(*)}^T \mathbf{x}_{i(*)} + \varepsilon_i \quad (4.1)$$

を、真のモデルと仮定する。いま、 $d$  個の説明変数から  $\mathbf{x}_{i(c)} = (1, x_{i(1)}, \dots, x_{i(k)})^T$  を選択し、真の密度関数を近似する候補モデルとして、パラメータ  $\theta_{(c)} = (\beta_{(c)}^T, \sigma_{(c)}^2)^T$  をもつ正規分布  $N(y_i|\mathbf{x}_{i(c)}; \beta_{(c)}^T \mathbf{x}_{i(c)}, \sigma_{(c)}^2)$  の密度関数  $f(y_i|\mathbf{x}_{i(c)}; \theta_{(c)})$  を考える。  $\alpha > 0$  における  $C_\alpha(g(\varepsilon); f(y|\mathbf{x}_{(c)}; \theta_{(c)}))$  を

$$\begin{aligned} & C_\alpha(g(\varepsilon); f(y|\mathbf{x}_{(c)}; \theta_{(c)})) \\ &= \frac{1}{\sqrt{(2\pi\sigma_{(c)}^2)^\alpha(1+\alpha)}} - \left(1 + \frac{1}{\alpha}\right) n^{-1} \sum_{i=1}^n \int_{\mathbf{R}} f(y_i|\mathbf{x}_{i(c)}; \theta_{(c)})^\alpha g(\varepsilon_i) d\varepsilon_i \end{aligned} \quad (4.2)$$

と定義する。そのとき、

$$\theta_{\alpha(c)} = \operatorname{argmin}_{\theta_{(c)} \in \Theta} C_\alpha(g(\varepsilon); f(y|\mathbf{x}_{(c)}; \theta_{(c)}))$$

とおき、この MDPDE を (3.5) 式から

$$\hat{\theta}_{\alpha(c)} = \operatorname{argmin}_{\theta_{(c)} \in \Theta} \ell_\alpha(\theta_{(c)}; f(y|\mathbf{x}_{(c)}; \theta_{(c)}))$$

によって与える。MDPDE について、Kobayashi et al. (2011) の定理 4.1 と同様に、次の定理が成り立つ。

定理 1. もしも  $\mathbf{x}_{i(c)}$  が真のモデルの説明変数  $x_{i(1*)}, \dots, x_{i(k_*)}$  をすべて含むならば、 $\tau \rightarrow 0$ ,  $\delta \rightarrow 0$  のもとで

$$D_\alpha(g(\varepsilon); f(y|\mathbf{x}_{(c)}; \theta_{\alpha(c)})) = O(\tau^2) + O(\tau^{1+\alpha}) + O(\delta^{1+\alpha})$$

が成り立つ。

注意 1. 真のモデルの説明変数  $x_{i(1*)}, \dots, x_{i(k_{**})}$  の中で  $x_{i(c)}$  に含まれないものが存在するとき、次を仮定をする.

$$\beta_{\alpha(c)}^T x_{i(c)} \neq \beta_{(*)}^T x_{i(*)}, \quad i = 1, \dots, n$$

そのとき,  $D_\alpha(g(\varepsilon); f(y|x_{(c)}; \theta_{\alpha(c)})) > 0$  となる.

いま,

$$\ell_\alpha(\theta; f(y_i|x_i; \theta)) = \frac{1}{\sqrt{(2\pi\sigma^2)^\alpha(1+\alpha)}} - \left(1 + \frac{1}{\alpha}\right) f(y_i|x_i; \theta)^\alpha$$

とおく. それに伴い,

$$\begin{aligned} I_\alpha(\theta) &= n^{-1} \sum_{i=1}^n \int_R \frac{\partial \ell_\alpha(\theta; f(y_i|x_i; \theta))}{\partial \theta} \frac{\partial \ell_\alpha(\theta; f(y_i|x_i; \theta))}{\partial \theta^T} dG(\varepsilon_i), \\ J_\alpha(\theta) &= n^{-1} \sum_{i=1}^n \int_R \frac{\partial^2 \ell_\alpha(\theta; f(y_i|x_i; \theta))}{\partial \theta \partial \theta^T} dG(\varepsilon_i) \end{aligned}$$

を定義する. そのとき, Kobayashi et al. (2011) の定理 4.2 と同様に, 次の定理が成り立つ.

定理 2.  $\tau = o(n^{-1/2})$  と仮定する. もしも  $x_{i(c)}$  が真のモデルの説明変数  $x_{i(1*)}, \dots, x_{i(k_{**})}$  をすべて含むならば,  $\tau^{1+\alpha} = o(n^{-1})$ ,  $\delta^{1+\alpha} = o(n^{-1})$  のもとで

$$D_\alpha(g(\varepsilon); f(y|x_{(c)}; \hat{\theta}_{\alpha(c)})) = (2n)^{-1} \text{tr}\{I_\alpha(\theta_{\alpha(c)})J_\alpha(\theta_{\alpha(c)})^{-1}\} + o(n^{-1})$$

が成り立つ.

2 節で見たように, MDPDE は  $\alpha$  の値を上手く選択することによって, MLE における異常値の影響を改善した. モデル選択においても, DPD に基づく情報量規準が, Akaike (1973) の AIC や Takeuchi (1976) の TIC などの K-L 情報量に基づく情報量規準を, 異常値の影響について改善することが期待される. いま,

$$\begin{aligned} \hat{I}_\alpha(\theta) &= n^{-1} \sum_{i=1}^n \frac{\partial \ell_\alpha(\theta; f(y_i|x_i; \theta))}{\partial \theta} \frac{\partial \ell_\alpha(\theta; f(y_i|x_i; \theta))}{\partial \theta^T}, \\ \hat{J}_\alpha(\theta) &= n^{-1} \sum_{i=1}^n \frac{\partial^2 \ell_\alpha(\theta; f(y_i|x_i; \theta))}{\partial \theta \partial \theta^T} \end{aligned}$$

とおき,  $E_G(\text{tr}\{\hat{I}_\alpha(\hat{\theta}_{\alpha(c)})\hat{J}_\alpha(\hat{\theta}_{\alpha(c)})^{-1}\}) = \text{tr}\{I_\alpha(\theta_{\alpha(c)})J_\alpha(\theta_{\alpha(c)})^{-1}\} + o(1)$ ,  $n \rightarrow \infty$  を仮定する.



[DPD 最小化モデル選択基準] 次で定義される  $IC_{\alpha(c)}$  が最小のモデルを、予測の意味で潜在分布に近いモデルとして選択する。

$$IC_{\alpha(c)} = 2n\ell_{\alpha}(\hat{\theta}_{\alpha(c)}; f(y|\mathbf{x}_{(c)}; \hat{\theta}_{\alpha(c)})) + 2\text{tr}\{\hat{\mathbf{I}}_{\alpha}(\hat{\theta}_{\alpha(c)})\hat{\mathbf{J}}_{\alpha}(\hat{\theta}_{\alpha(c)})^{-1}\} \quad (4.3)$$

もしも、 $k+1$  個の説明変数をもつ  $\mathbf{x}_{i(c)}$  が真のモデルの説明変数をすべて含み、 $k+2$  個の説明変数をもつ  $\mathbf{x}_{i(c')}$  が  $\mathbf{x}_{i(c)}$  の説明変数をすべて含むとき、次の条件を仮定する。

$$0 < \text{tr}\{\mathbf{I}_{\alpha}(\theta_{\alpha(c)})\mathbf{J}_{\alpha}(\theta_{\alpha(c)})^{-1}\} < \text{tr}\{\mathbf{I}_{\alpha}(\theta_{\alpha(c')})\mathbf{J}_{\alpha}(\theta_{\alpha(c')})^{-1}\}$$

そのとき、Kobayashi et al. (2011) の定理 4.3 と同様に、次の定理を得る。

定理 3.  $\tau^2 = o(n^{-1})$ ,  $\tau^{1+\alpha} = o(n^{-1})$ ,  $\delta^{1+\alpha} = o(n^{-1})$  を仮定する。そのとき、真のモデルの  $IC_{\alpha(*)}$  の期待値  $E_G[IC_{\alpha(*)}]$  は、候補モデル  $N(y_i|\mathbf{x}_{i(c)}; \beta_{(c)}^T \mathbf{x}_{i(c)}, \sigma_{(c)}^2)$  の  $IC_{\alpha(c)}$  の期待値  $E_G[IC_{\alpha(c)}]$  を最小にする。

注意 2.  $\alpha$  の値は、定理 3 の収束条件を満たす範囲で選択する。真のモデルの情報量基準  $IC_{\alpha(*)}$  が平均的に最小値を与えることから、DPD 最小化モデル選択基準によって真のモデルを平均的に選択することができる。

## 5. シミュレーション

4 節で提案した情報量基準  $IC_{\alpha(c)}$  の性能を、シミュレーション実験で検証する。次のような線形回帰モデルを考えた。

$$y_i = \beta_0 + \sum_{j=1}^6 \beta_j x_{ij} + \varepsilon_i$$

ここで、 $x_{ij}, j = 1, \dots, 6$  は、一様分布  $U(0, 2)$  から独立に発生させた。誤差変数  $\varepsilon_i$  は、異常値が混入した以下の分布を考えた。

$$g(\varepsilon) = 0.95\phi(\varepsilon|0, 1) + 0.05\psi(\varepsilon|4 < \varepsilon < 8)$$

真のモデルに、次の 2 つのモデルを考えた。

$$(I) \quad \beta_0 = 1, \beta_1 = -1, \beta_2 = 1, \beta_3 = \dots = \beta_6 = 0$$

$$(II) \quad \beta_0 = 1, \beta_1 = -1, \beta_2 = -0.5, \beta_3 = 0.5, \beta_4 = 1, \beta_5 = \beta_6 = 0$$

候補モデル  $N(y_i|\mathbf{x}_{i(c)}; \beta_{(c)}^T \mathbf{x}_{i(c)}, \sigma_{(c)}^2)$  については、 $\mathbf{x}_{i(c)} (= \mathbf{x}_{ik(c)}) = (1, x_{i1}, \dots, x_{ik})^T$ ,  $k = 1, \dots, 6$  の 6 個のモデルを考えた。真のモデルは、(I) の場合は  $\mathbf{x}_{i2(c)}$ 、(II) の場合は  $\mathbf{x}_{i4(c)}$  で与えられることに注意する。いま、真の分布から独立にデータを 60 ( $= n$ ) 個発生させ、潜在分布の推定を、MLE ( $\alpha = 0$ ) と MDPDE ( $\alpha = 0.2, 0.4, 0.6$ ) で

行った。構築した6個のモデルの中から、TIC ( $\alpha = 0$ ) と  $IC_{\alpha(c)}$  ( $\alpha = 0.2, 0.4, 0.6$ ) を用いて、各  $\alpha$  について情報量規準を最小にするモデルを選択した。この実験を独立に1000回繰り返し、TIC ( $\alpha = 0$ ) と  $IC_{\alpha(c)}$  ( $\alpha = 0.2, 0.4, 0.6$ ) による各モデルの選択回数を、表3に纏めた。

表3. TIC ( $\alpha = 0$ ) と  $IC_{\beta(c)}$  ( $\alpha = 0.2, 0.4, 0.6$ ) による  
各モデルの選択回数 (実験回数: 1000 回)

| 真のモデルが (I) の場合       |     |     |     |    |    |    |
|----------------------|-----|-----|-----|----|----|----|
| $\alpha \setminus k$ | 1   | 2   | 3   | 4  | 5  | 6  |
| 0                    | 160 | 494 | 155 | 74 | 63 | 54 |
| 0.2                  | 16  | 631 | 179 | 85 | 66 | 23 |
| 0.4                  | 11  | 628 | 192 | 91 | 53 | 25 |
| 0.6                  | 24  | 645 | 203 | 79 | 37 | 12 |

| 真のモデルが (II) の場合      |     |    |    |     |     |     |
|----------------------|-----|----|----|-----|-----|-----|
| $\alpha \setminus k$ | 1   | 2  | 3  | 4   | 5   | 6   |
| 0                    | 142 | 71 | 35 | 483 | 149 | 120 |
| 0.2                  | 33  | 15 | 4  | 713 | 153 | 82  |
| 0.4                  | 18  | 12 | 3  | 733 | 157 | 77  |
| 0.6                  | 37  | 28 | 17 | 690 | 172 | 56  |

MLEに基づくTICは異常値に敏感で、真のモデルを有効に選択していないことが見て取れた。一方、MDPDEに基づく $IC_{\beta(c)}$ は、真のモデルをTICよりも高い確率で選択していることが見て取れた。設定を変えた実験においても同様の結果が得られ、DPD最小化モデル選択基準 $IC_{\beta(c)}$ が、異常値に対して頑健なモデル選択を実現していることが確認できた。

謝辞 本研究は、科学研究費補助金 基盤研究 (B) 22300094 研究代表者: 青嶋 誠「高次元データの理論と方法論の総合的研究」から、研究助成を受けています。

## 参考文献

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. *2nd Inter. Symp. on Information Theory* (Petrov, B. N. and Csaki, F., eds.), Akademiai Kiado, 267–281 (Reproduced in *Breakthroughs in Statistics*, 1 (Kotz, S. and Johnson, N. L., eds.), Springer-Verlag (1992)).
- Basu, A., Harris, I. R., Hjort, N. L. and Jones, M. C. (1998). Robust and efficient estimation by minimising a density power divergence. *Biometrika* **85**, 549–559.

- Basu, S., Basu, A. and Jones, M. C. (2006). Robust and efficient parametric estimation for censored survival data. *Ann. Inst. Statist. Math.* **58**, 341-355.
- Durio, A. and Isaia, E.D. (2011). The minimum density power divergence approach in building robust regression models. *Informatica* **22**, 43-56.
- Fujisawa, H. and Eguchi, S. (2006). Robust estimation in the normal mixture model. *J. Statist. Plan. Infer.* **136**, 3989-4011.
- Fujisawa, H. and Eguchi, S. (2008). Robust parameter estimation with a small bias against heavy contamination. *J. Multivariate Anal.* **99**, 2053-2081.
- Hjort, N. L. (1994). Minimum L2 and robust Kullback-Leibler estimation. *Proceedings of the 12th Prague Conference on Information Theory, Statistical Decision Functions and Random Process* (P. Lachout and J. Visek, eds.), 102-105. Prague: Academy of Sciences of the Czech Republic.
- Hong, C. and Kim, Y. (2001). Automatic selection of the tuning parameter in the minimum density power divergence estimation. *Journal of the Korean Statistical Association.* **30**, 453-465.
- Jones, M. C., Hjort, N. L., Harris, I. R. and Basu, A. (2001). A comparison of related density-based minimum divergence estimators. *Biometrika* **88**, 865-873.
- Kobayashi, Y., Yata, K. and Aoshima, M. (2011). Robust and accelerated model selection for clustering in a contaminated mixture model, submitted.
- Scott, D. W. (2001). Parametric statistical modeling by minimum integrated square error. *Technometrics* **43**, 274-285.
- Takeuchi, K. (1976). Distribution of information statistics and criteria for adequacy of models. *Mathematical Science* **153** 12-18 (in Japanese).
- Warwick, J. (2005). A data-based method for selecting tuning parameters in minimum distance estimators. *Computational Statistics and Data Analysis* **48**, 571-585.
- Warwick, J. and Jones, M. C. (2005). Choosing a robustness tuning parameter. *Journal of Statistical Computation and Simulation* **75**, 581-588.